# DATA MINING & DISCOVERY SQL ASSIGNMENT

### **SUBMITTED BY:**

NAME: Saranya Bala Subramaniam

STUDENT ID: 22033095

#### PYTHON SCRIPT THAT GENERATED THE DATABASE:

```
import csv
import random
from faker import Faker
import sqlite3
import os
fake = Faker()
# Function to generate random ordinal data
def generate_ordinal_data():
  ratings = ["Low", "Medium", "High"]
  return random.choice(ratings)
# Function to generate random nominal data
def generate nominal data():
  categories = ["Electronics", "Clothing", "Home Appliances"]
  return random.choice(categories)
# Function to generate random ratio data
def generate_ratio_data():
  return round(random.uniform(10.0, 500.0), 2)
# Function to generate random interval data
def generate interval data():
  return random.randint(1, 10)
# Set the current working directory to the script's directory
script_directory = os.path.dirname(os.path.realpath(__file__))
os.chdir(script directory)
```

```
# Generate Products data
```

```
products_data = []
for product_id in range(1, 1001):
  category = generate nominal data()
  price = generate ratio data()
  stock quantity = random.randint(20, 100)
  product rating = generate ordinal data()
  customer_id = fake.uuid4() # Adding customer_id
    product = {
    "product_id": product_id,
    "product_name": fake.word(),
    "category": category,
    "price": price,
    "stock_quantity": stock_quantity,
    "product_rating": product_rating,
    "customer_id": customer_id,
 }
  products data.append(product)
# Generate Sales data
sales data = []
for sale_id in range(1, 1001):
  customer name = fake.name()
  sale date = fake.date this year()
  customer_id = fake.uuid4() # Adding customer_id
    sale = {
    "sale id": sale id,
    "customer_id": customer_id,
    "customer_name": customer_name,
    "sale_date": sale_date,
 }
```

```
sales data.append(sale)
# Generate Sales Details data (Child Table)
sales details data = []
for sale id in range(1, 1001):
  product id = random.randint(1, 1000)
  quantity sold = random.randint(1, 10)
    sales details = {
    "sale id": sale id,
    "product id": product id,
    "quantity_sold": quantity_sold,
  }
  sales details data.append(sales details)
# Define CSV file names with full paths
products_csv_file = os.path.join(script_directory, "products_data.csv")
sales_csv_file = os.path.join(script_directory, "sales_data.csv")
sales details csv file = os.path.join(script directory, "sales details data.csv")
# Write data to CSV files
with open(products csv file, mode="w", newline="") as file:
  fieldnames = ["product_id", "product_name", "category", "price", "stock_quantity",
"product_rating", "customer_id"]
  writer = csv.DictWriter(file, fieldnames=fieldnames)
  writer.writeheader()
  writer.writerows(products data)
with open(sales_csv_file, mode="w", newline="") as file:
  fieldnames = ["sale_id", "customer_id", "customer_name", "sale_date"]
  writer = csv.DictWriter(file, fieldnames=fieldnames)
  writer.writeheader()
  writer.writerows(sales data)
```

```
with open(sales details csv file, mode="w", newline="") as file:
  fieldnames = ["sale_id", "product_id", "quantity_sold"]
  writer = csv.DictWriter(file, fieldnames=fieldnames)
  writer.writeheader()
  writer.writerows(sales details data)
# Function to create SQLite database and import data from CSV files
def create sales database(products csv, sales csv, sales details csv, database name):
  conn = sqlite3.connect(database name)
  cursor = conn.cursor()
  # Create 'products' table if it does not exist
  cursor.execute("SELECT name FROM sqlite master WHERE type='table' AND
name='products'")
  if not cursor.fetchone():
    with open(products csv, 'r') as file:
      reader = csv.reader(file)
      header = next(reader)
      columns = ', '.join(header)
      cursor.execute(f'CREATE TABLE products ({columns})')
      cursor.executemany(f'INSERT INTO products VALUES ({", ".join(["?"] * len(header))})',
reader)
  # Create 'sales' table if it does not exist
  cursor.execute("SELECT name FROM sglite master WHERE type='table' AND
name='sales'")
  if not cursor.fetchone():
    with open(sales csv, 'r') as file:
      reader = csv.reader(file)
      header = next(reader)
      columns = ', '.join(header)
      cursor.execute(f'CREATE TABLE sales ({columns})')
```

```
cursor.executemany(f'INSERT INTO sales VALUES ({", ".join(["?"] * len(header))})',
reader)
  # Create 'sales details' table if it does not exist
  cursor.execute("SELECT name FROM sqlite_master WHERE type='table' AND
name='sales_details'")
  if not cursor.fetchone():
    with open(sales_details_csv, 'r') as file:
      reader = csv.reader(file)
      header = next(reader)
      columns = ', '.join(header)
      cursor.execute(f'CREATE TABLE sales_details ({columns})')
      cursor.executemany(f'INSERT INTO sales_details VALUES ({", ".join(["?"] *
len(header))})', reader)
  # Commit changes and close connection
  conn.commit()
  conn.close()
# Create the 'sales_database.db' and import data
create_sales_database(products_csv_file, sales_csv_file, sales_details_csv_file,
'sales database.db')
```

#### 1.Data Generation:

- Using the Faker package and random functions, the given Python script creates fake data for a sales database.
- The generated data includes information about products, sales, and sales details.
- Product categories, prices, and stock quantities are representative of a fictional sales scenario. Customer names, product names, and sale dates are randomly generated.
- The synthetic data generated in this script can be used for testing and development purposes, simulating a sales database with related tables for products and sales details.

Below is an explanation of the data generation process:

1. Ordinal Data (Product Ratings):

```
# Function to generate random ordinal data
def generate_ordinal_data():
    ratings = ["Low", "Medium", "High"]
    return random.choice(ratings)
```

2. Nominal Data (Product Categories):

```
# Function to generate random nominal data
def generate_nominal_data():
    categories = ["Electronics", "Clothing", "Home Appliances"]
    return random.choice(categories)
```

3. Ratio Data (Product Price):

```
# Function to generate random ratio data
def generate_ratio_data():
    return round(random.uniform(10.0, 500.0), 2)
```

4. Interval Data (Stock Quantity):

```
# Function to generate random interval data
def generate_interval_data():
    return random.randint(1, 10)
```

- 5. UUID Generation (Customer ID):
  - Customer IDs are generated using the fake.uuid4() function from the Faker library.
- 6. Date Generation (Sale Date):
- Sale dates are generated using the <a href="fake.date\_this\_year()">fake.date\_this\_year()</a> function from the Faker library.
- 7. Random Product Names and Customer Names:

- Product names and customer names are generated using the fake.word() and fake.name() functions, respectively.

#### CSV File Creation:

- Three CSV files are created: "products\_data.csv," "sales\_data.csv," and "sales details data.csv."
- Data for products, sales, and sales details are written to these files using the `csv.DictWriter` module.

#### **SQLite Database Creation:**

- A SQLite database named "sales database.db" is created.
- The script defines a function, `create\_sales\_database`, to create tables based on the CSV file structures and import data into them.
- The function is called for each CSV file to create tables for products, sales, and sales details.

#### 2.THE SCHEMA OF THE DATABASE:

The schema establishes relationships between the tables using foreign keys and ensures data integrity through primary keys.

#### Table 1- products:

```
CREATE TABLE "products" (

"product_id" INTEGER NOT NULL,

"product_name" TEXT,

"category" TEXT,

"price" REAL,

"stock_quantity" INTEGER,

"product_rating" TEXT,

"customer_id" TEXT NOT NULL,

PRIMARY KEY("product_id"),

FOREIGN KEY("customer_id") REFERENCES "sales"("customer_id")

);
```

'product\_id' is the primary key, and 'customer\_id' is a foreign key referencing the 'customer\_id' in the 'sales' table.

- product\_id: An integer representing the unique identifier for each product. It is marked as NOT NULL, indicating it must have a value, and is set as the primary key for this table.
- product\_name: A text field storing the name of the product.
- category: A text field indicating the category of the product.
- o price: A real number representing the price of the product.
- stock\_quantity: An integer representing the quantity of the product in stock.
- product\_rating: A text field representing the rating of the product.
- o <u>customer\_id</u>: A text field representing the customer ID to which this product is associated. It is marked as NOT NULL, indicating it must have a value, and it is a foreign key referencing the customer\_id in the sales table.

#### Table 2- sales:

```
CREATE TABLE "sales" (

"sale_id" INTEGER NOT NULL,

"customer_id" TEXT NOT NULL,

"customer_name" TEXT,

"sale_date" DATE,

PRIMARY KEY("sale_id","customer_id"),

FOREIGN KEY("customer_id") REFERENCES "products"("customer_id")

);
```

the primary key is a composite key consisting of 'sale\_id' and 'customer\_id'. This means that the combination of sale\_id and customer\_id must be unique. The 'customer\_id' is a foreign key referencing the 'customer\_id' in the 'products' table.

- o sale\_id: An integer representing the unique identifier for each sale. It is marked as NOT NULL, indicating it must have a value, and it is part of the primary key for this table.
- o <u>customer\_id</u>: A text field representing the customer ID associated with this sale. It is marked as NOT NULL, indicating it must have a value, and it is a foreign key referencing the customer\_id in the products table.
- customer\_name: A text field storing the name of the customer.
- sale date: A date field indicating the date of the sale.

#### Table 3- sales details:

```
CREATE TABLE "sales_details" (

"sale_id" INTEGER NOT NULL,

"product_id" INTEGER NOT NULL,

"quantity_sold" INTEGER,

PRIMARY KEY("product_id", "sale_id"),

FOREIGN KEY("product_id") REFERENCES "products"("product_id"),

FOREIGN KEY("sale_id") REFERENCES "sales"("sale_id")

);
```

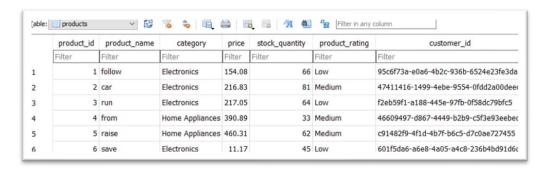
the primary key is a composite key consisting of `product\_id` and `sale\_id`. The `product\_id` is a foreign key referencing the `product\_id` in the `products` table, and the `sale\_id` is a foreign key referencing the `sale\_id` in the `sales` table.

- sale\_id: An integer representing the unique identifier for each sale detail. It is marked
  as NOT NULL, indicating it must have a value, and it is part of the primary key for this
  table. It is also a foreign key referencing the sale\_id in the sales table.
- product\_id: An integer representing the unique identifier for each product detail. It is marked as NOT NULL, indicating it must have a value, and it is part of the primary key for this table. It is also a foreign key referencing the product\_id in the products table.
- quantity\_sold: An integer representing the quantity of the product sold in the specific sale.

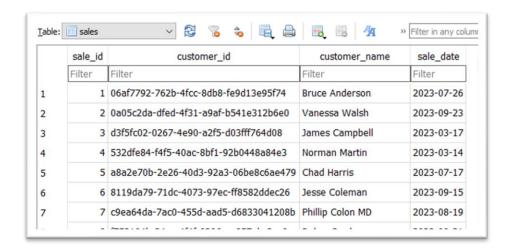
## 3. JUSTIFICATION FOR ANY SEPARATE TABLE AND ETHICAL DISCUSSION:

By minimising redundancy and adhering to normalisation standards, this separation protects data integrity through relationships.

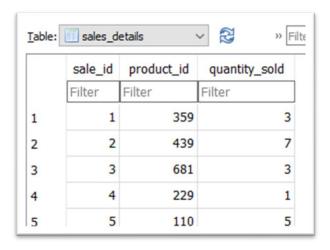
**Products Table:** Contains information about individual products such as product\_id, product\_name, category, price, stock\_quantity, product\_rating, and customer\_id.



**Sales Table:** Stores information about sales transactions, including sale\_id, customer\_id, customer\_name, and sale\_date.



**Sales Details Table:** Captures details about each sale, linking sale\_id to product\_id and recording quantity\_sold.



By minimising redundancy and adhering to normalisation standards, this separation protects data integrity through relationships.

#### Flexibility and Maintenance:

**Products**: It is simple to edit product details without compromising sales data.

**Sales**: Individual product details are unaffected by changes to client information or sale dates.

Sales details: Allows for changes to amounts sold without affecting essential product or sales data.

The system is more versatile and easier to maintain because to its modular architecture.

#### **Privacy and Security:**

- Products: Ethical considerations include safeguarding product data.
- Sales: Protects customer-related information.
- Sales Details: Ensures transaction details are secure.

Robust security measures are vital to prevent unauthorized access and protect sensitive data.

#### **Consent and Transparency:**

- Products: Users should be aware of data collection related to products.
- Sales: Transparency in sales transactions, with clear communication about data usage.
- o **Sales Details**: Users need to understand how transaction details are recorded.

Transparency builds trust and helps users make informed decisions.

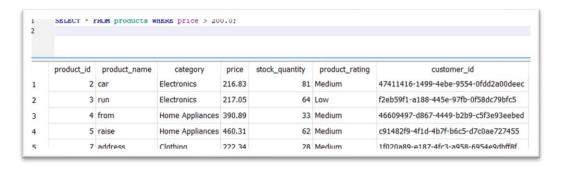
#### Data Accuracy and Fairness:

- Products: Ensuring accurate representation of products.
- Sales: Accuracy in recording sales transactions.
- Sales Details: Precision in capturing quantities sold.

Ethical responsibility involves maintaining accurate and fair data representation.

## 4. Example queries of your database including joins and selections, demonstrating different data types.

Query 1 - Selecting Products with Price Greater Than \$200:



**Query 2 -** Joining Sales and Sales Details to Get Customer Names and Quantity Sold:

1	SELECT s.sale_id, s.customer_name, sd.quantity_sold				
2	FROM sales s				
3	JOIN sales_details sd ON s.sale_id = sd.sale_id;				
4		_	_		
	sale_id	customer_name	quantity_sold		
1	1	Bruce Anderson	3		
2	2	Vanessa Walsh	7		
3	3	James Campbell	3		
4	4	Norman Martin	1		
5	5	Chad Harris	5		
6	6	Jesse Coleman	4		

**Query 3** - Aggregating Sales Details to Get Total Quantity Sold for Each Product:

	FROM sales		
	GROUP BY p	roduct_id;	
_	product_id	total_quantity_sold	
	2	8	
	3	4	
	5	3	
	7	8	
	9	1	